# MALT : A MALloc Tracker

## A memory profiling tool

# THE QUESTION

Software R&D: Observing Memory - MALT,
Sébastien Valat

- We have **profiling tool** for **timing** (eg. Valgrind or vtune)

- But for **memory usage** ?

- Memory can be an issue :
  - Failed to run (or swap) due to **lack of memory resource**.
  - **Performance impact** of memory management functions

- Three main questions :
  - How to reduce **memory footprint** ?
  - How to improve overhead of **memory management** ?
  - How to improve **memory usage** ?

Software R&D: Observing Memory - MALT,
Sébastien Valat

**Exascale ∞**
*computing research*

- We want to help searching :
  - **Where** memory is allocated.
  - **Properties** of allocated chunks.
  - **Bad** allocation **patterns** for performance.
  - **Leaks**
  - **Global variables** (TLS)

```
1    __thread Int gblVar[SIZE];
2    int * func(int size)
3    {
4              child_func_with_allocs();
5              void * ptr = new char[size];
6              double* ret = new double[size*size*size];
7              for (.....)
9              {
10                        double* buffer = new double[size];
11                        //short and quick do stuff
12                        delete [] buffer;
13              }
14             return ret;
15   }
```
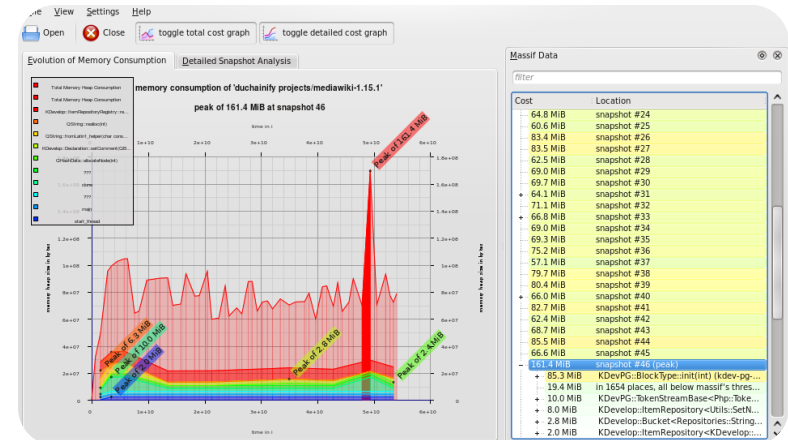
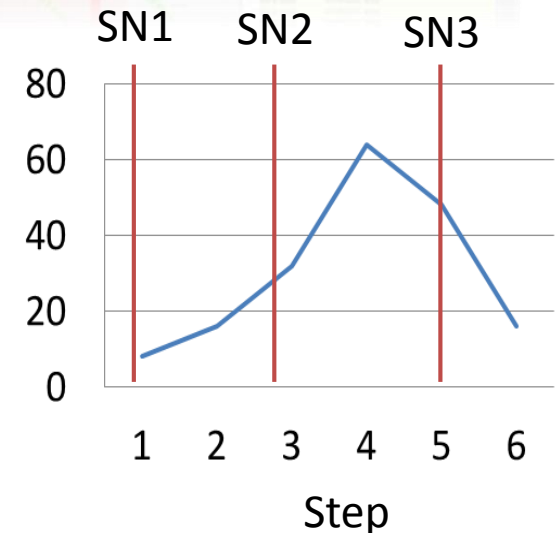| Global variables and TLS |
| Indirect allocations |
| Leak |
| Might lead to swap for large size |
| Short life allocations |

Software R&D: Observing Memory – MALT, Sébastien Valat

# EXISTING TOOLS

Software R&D: Observing Memory - MALT,
Sébastien Valat

- **Valgrind - massif :**
  – Take **snapshots** over time.
  – **Link memory** size to **functions**
  – **Peak** might **not be captured**.
  – Might **miss short live** allocations
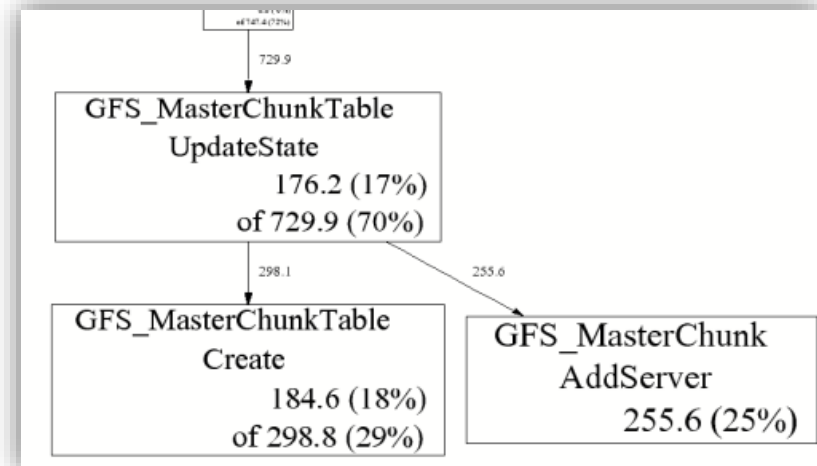  – GUI not adequate for large code
  – Slow, not parallel.

- **Valgrind - memcheck :**
  – Misuse of memory functions (malloc/new/free….)
  – Leak detection
  – Invalid accesses
  – Slow, not parallel.

Software R&D: Observing Memory - MALT, Sébastien Valat
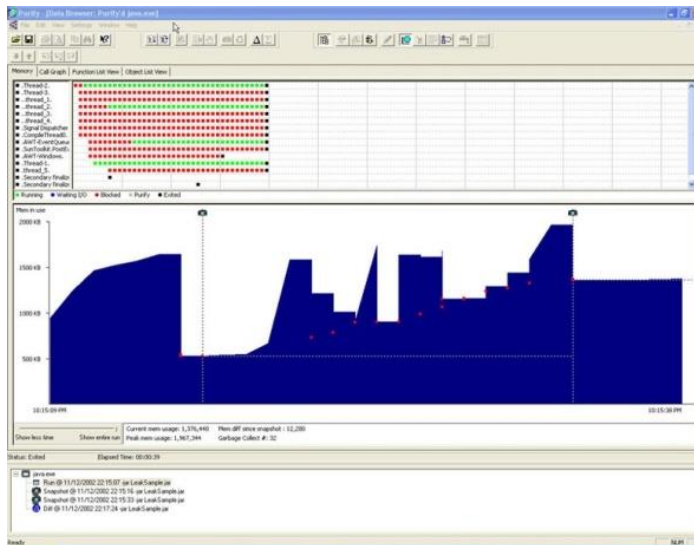
- **Google heap profiler (tcmalloc):**
  - Small overhead.
  - Only provide snapshots of allocated memory per stacks.
  - Peak might not be captured.
  - No binary or source instrumentation.
  - Output is not clear…
  - Lack of a real GUI to use it.
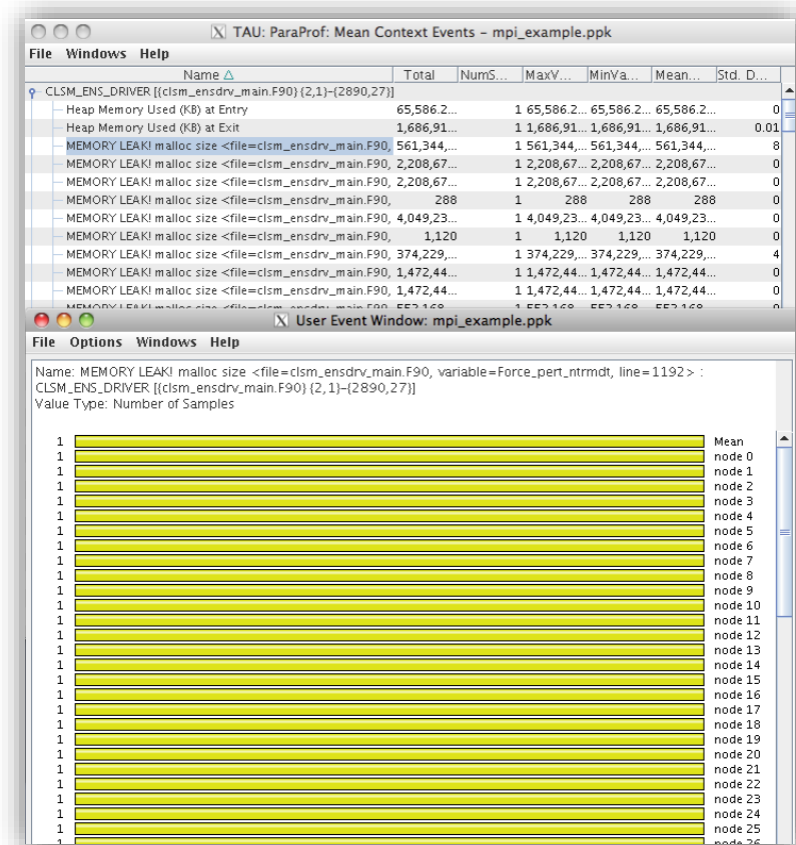


```
% pprof gfs_master profile.0100.heap
  255.6  24.7%  24.7%   255.6  24.7% GFS_MasterChunk::AddServer
  184.6  17.8%  42.5%   298.8  28.8% GFS_MasterChunkTable::Create
  176.2  17.0%  59.5%   729.9  70.5% GFS_MasterChunkTable::UpdateState
  169.8  16.4%  75.9%   169.8  16.4% PendingClone::PendingClone
   76.3   7.4%  83.3%    76.3   7.4% __default_alloc_template::_S_chunk_alloc
   49.5   4.8%  88.0%    49.5   4.8% hashtable::resize
```

Software R&D: Observing Memory - MALT, Sébastien Valat

- **IBM Purify++ / Parasoft Insure++**
  - Commercial
  - Leak detection, access checking, memory debugging tools.
  - Use binary or source instrumentation.
  - Windows / Redhat
  - Quite old GUI

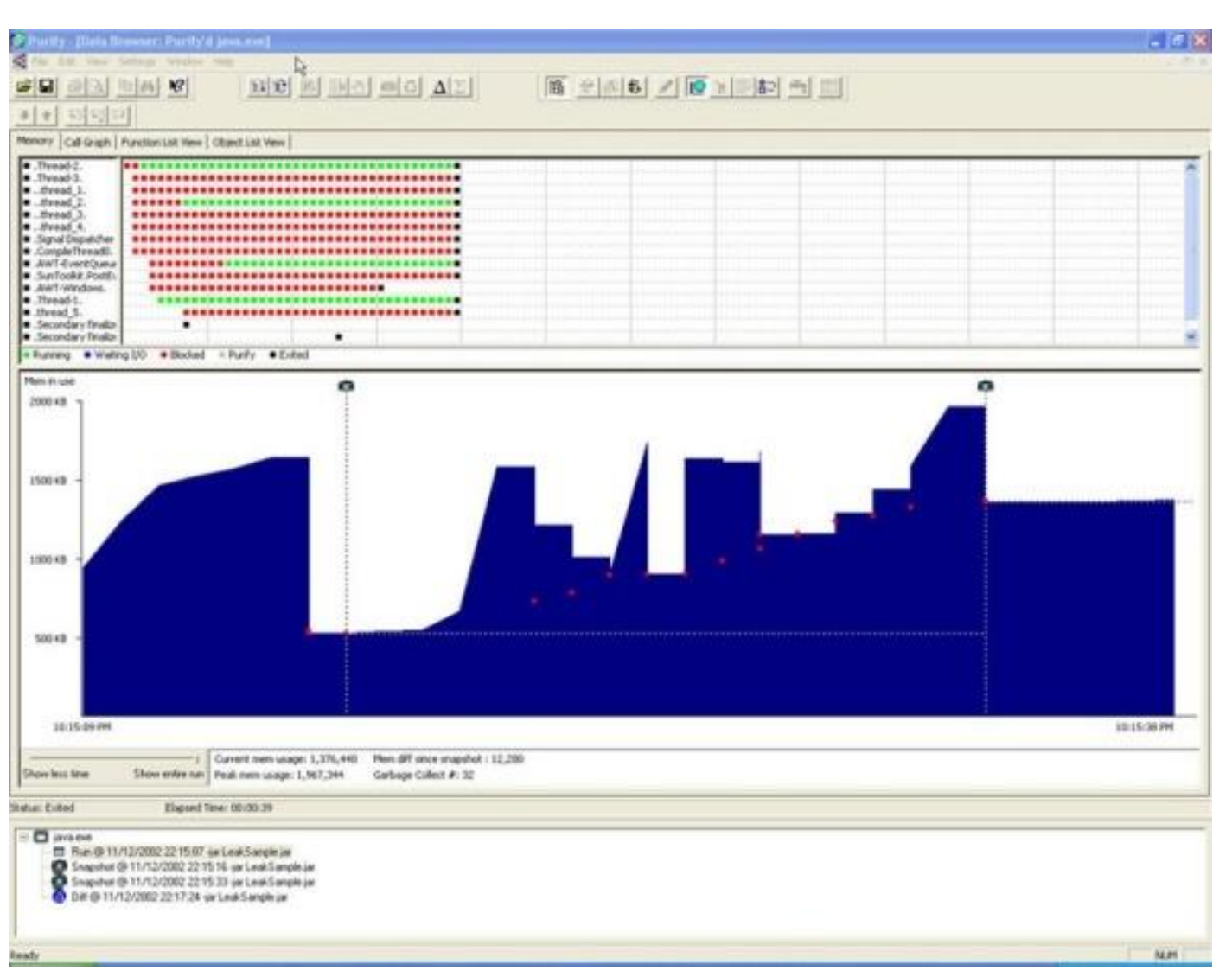

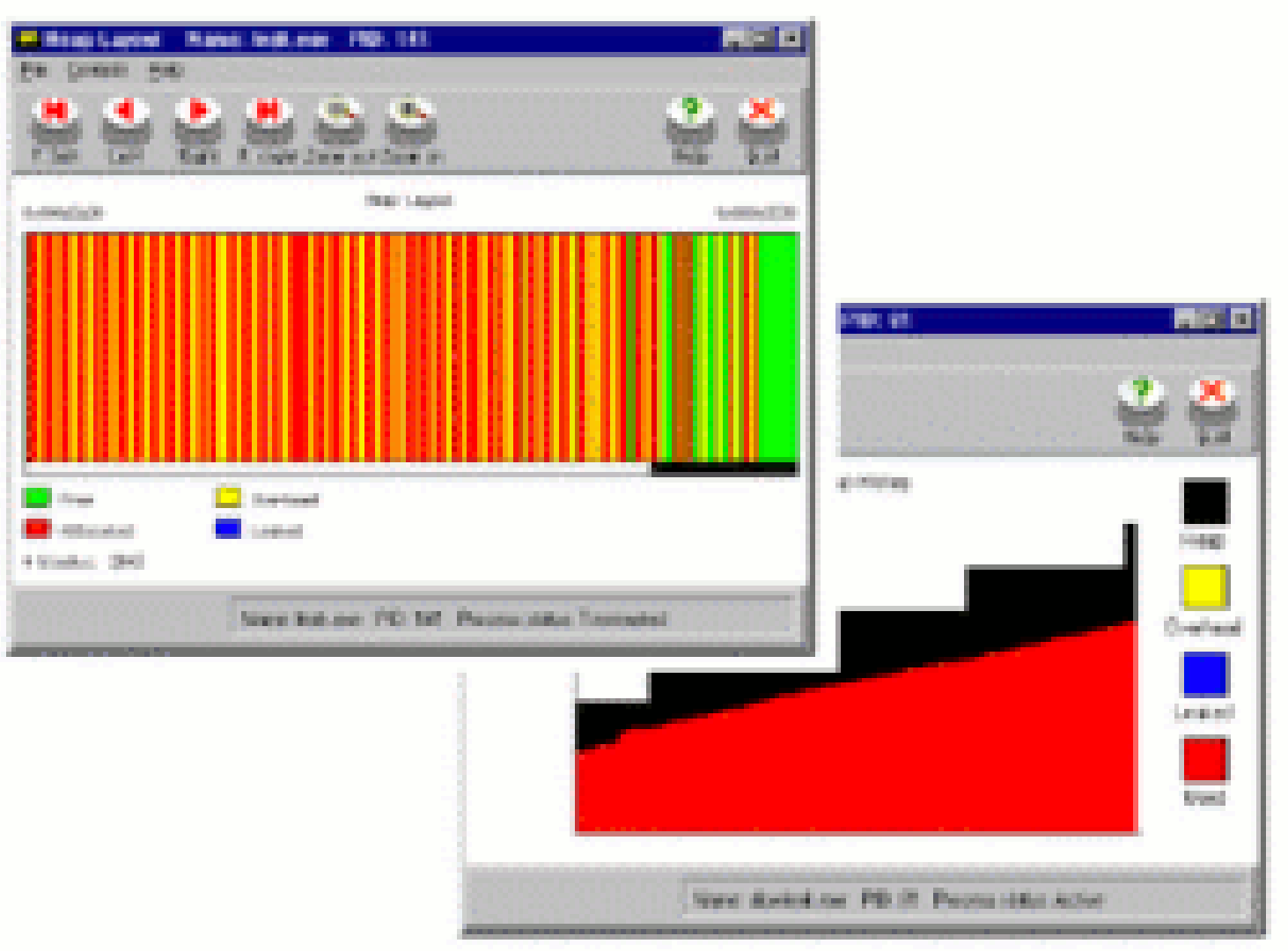

- **TAU memory profiler**
  - Provide profiles
  - Follow stacks
  - Track leaks
  - Parallel, done for HPC/MPI
  - Lack easy matching with sources

Software R&D: Observing Memory - MALT, Sébastien Valat

# What is good in kcachgrind



- List of **functions** with **exclusive/inclusive** costs



- Nice **call tree**



- **Annotated** sources

Software R&D: Observing Memory - MALT,
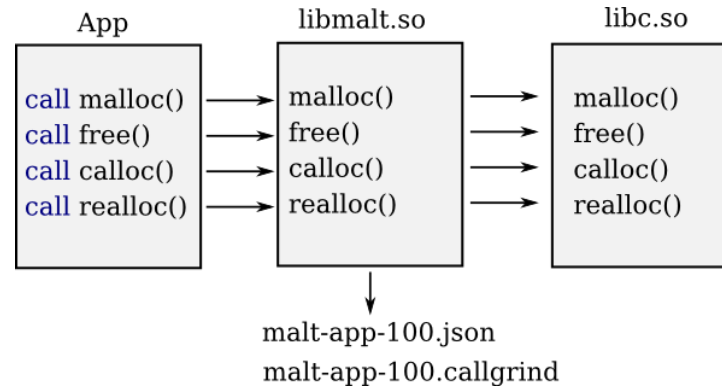Sébastien Valat

# WHERE I GO

- Same **approach** than **valgrind/kcachgind**

- **Mapped** allocations on **sources lines**

- For **memory resource usage** :
  – Memory **leaks** (malloc without free)
  – **Peak** and **total** allocated **memory**

- For **performance** :
  – Allocation **count**
  – Allocation **sizes** (min/mean/max)
  – Chunk **lifetime** (min/mean/max)

**Exascale ∞** computing research

- **Profile over time :**
  - Allocation **rate**
  - **Physical / Virtual / Requested** memory
  - **Stack size** for each **thread** (require function instrumentation)

- **Example on YALES2 with gfortran :**

Software R&D: Observing Memory - MALT, Sébastien Valat

- Use LD_PRELOAD to intercept malloc/free/…



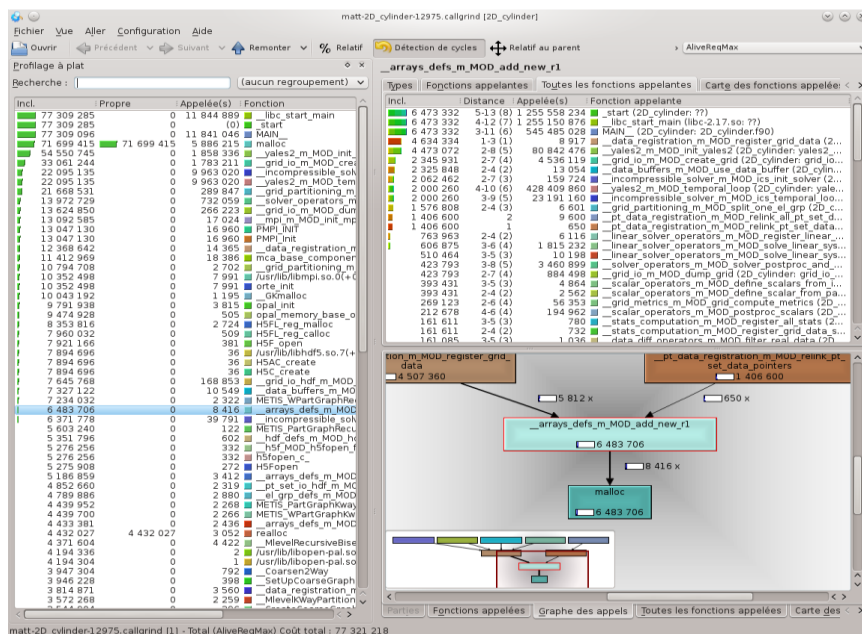- Profile allocations on call stacks

- Generate JSON output file

- Build profile so size is limited by call tree

Software R&D: Observing Memory - MALT, Sébastien Valat

- Two approach implemented : backtrace and instrumentation

- Backtrace (default) :
  - No change on binary or sources.
  - Function from glibc
  - Manage all dynamic libraries
  - No impact on compute
  - Slow on x86_64 due to use of libunwind.
  - Slow for large number of calls (~>10M)

- Instrumentation :
  - Need source recompilation (available) : *-finstrument-function*
  - Or tools for binary instrumentation : MAQAO / Pintool (experimental)
  - Impact performance of compute only functioans
  - Faster for really large number of calls to malloc
  - Only provide stacks for the instrumented sources/binaries

Software R&D: Observing Memory - MALT,
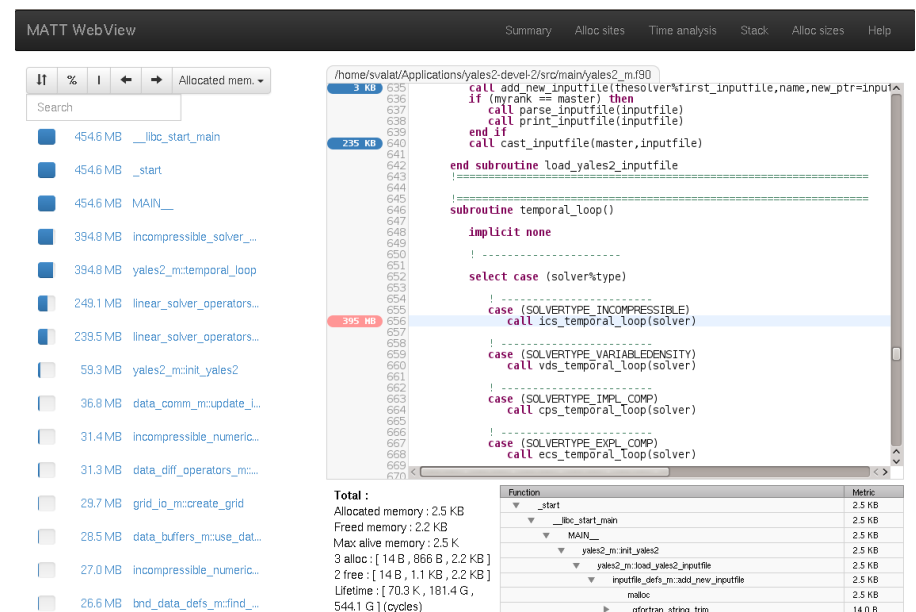Sébastien Valat

## Callgrind compatibiltiy

- Can use kcachgrind
- Might be usefull for some users, cannot provide all metrics.

## Own web view

- Get all metrics
- Web technology (NodeJS, D3JS, Jquery, AngularJS)
- Easier for remote usage
- Can be used for shared working

Software R&D: Observing Memory - MALT, Sébastien Valat

- Display **human readable** units
  - You prefer **15728640** of **15MB** ?
  - I want to **compare to what I expect**.

- Cannot handle **non sum cumulative metrics**
  - **Inclusive** costs **only rely** on **+ operator**
  - Some mem. metrics **requires max/min** (eg. local peaks, lifetime, sizes)

- No way to express **time profiles**

- No way to express **parameter distributions** (eg. sizes).

# SOME VIEWS

Software R&D: Observing Memory - MALT,
Sébastien Valat

| EXECUTION TIME | PHYSICAL MEMORY PEAK | ALLOCATION COUNT | AVAILABLE PHYSICAL MEMORY |
|---|---|---|---|
| 00:00:00.25 | 2.3 MB | 379 | 4.1 Gb |

## Run description

| Executable : | simple-case-finstr-linked |
|---|---|
| Commande : | ./simple-case-finstr-linked |
| Tool : | matt-0.0.0 |
| Host : | localhost |
| Date : | 2014-11-26 22:40 |
| Execution time : | 00:00:00.25 |
| Ticks frequency : | 1.8 GHz |

## Global statistics

Show all details   Show help

| Physical memory peak | 2.3 MB |
|---|---|
| Virtual memory peak | 103.7 MB |
| Requested memory peak | 2.0 KB |

- Provide a small summary
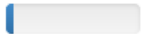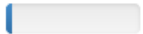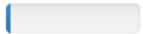- Provide some warnings

---

| Show all details | Show help | |
|---|---|---|
| Physical memory peak | | 66.7 MB |
| Virtual memory peak | | 158.1 MB |
| Requested memory peak | | 6.1 MB |
| Cumulated memory allocations | | 11.5 MB |
| Allocation count | | 172.2 K |
| Recycling ratio | | 1.9 |
| Leaked memory | | 743.7 KB |
| Largest stack | | 0 B |
| **Global variables** | | **10.0 MB** ⚠ |
| TLS variables | | 48 B |
| **Global variable count** | | **421.0 K** ⚠ |
| Peak allocation rate | | 37.8 MB/s |

- Summarize **top functions** for some metrics

- Points to check

- Examples on YALES2

## Alloc count

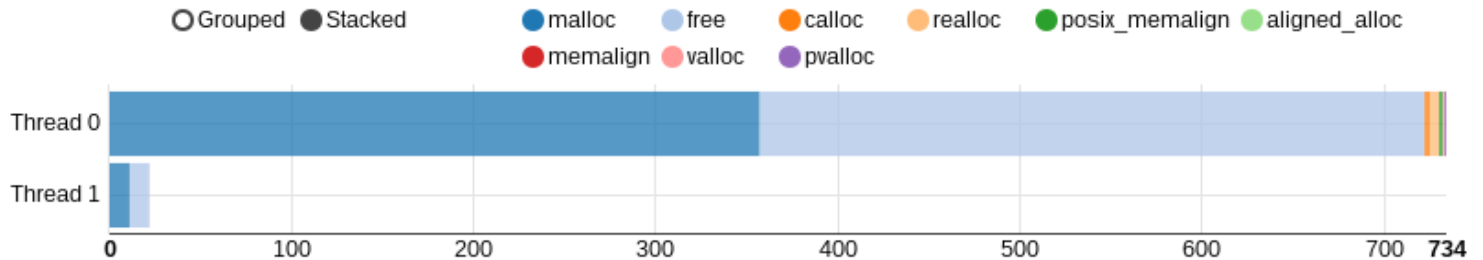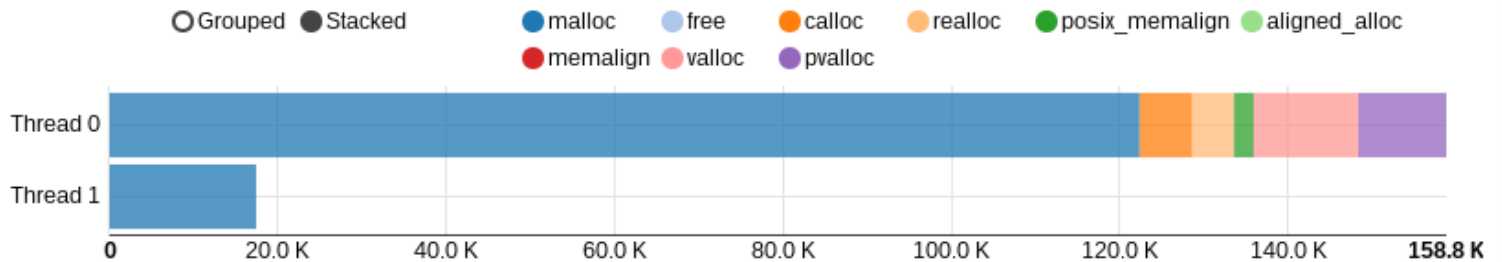| Ratio | Allocs | Function |
|---|---|---|
| | 911.9 K | data_comm_m::copy_int_comm_to_data |
| | 896.4 K | data_comm_m::copy_data_to_int_comm |
| | 853.2 K | data_comm_m::update_int_comm |
| | 484.9 K | sponge_layer_m::calc_sponge_layer_mask |
| | 296.0 K | incompressible_numerics_m::ics_diffuse_velocity_rk_4th |

## Allocated memory

| Ratio | Allocs | Function |
|---|---|---|
| | 202.4 MB | linear_solver_operators_m::solve_linear_system_deflated_pcg |
| | 26.6 MB | bnd_data_defs_m::find_bnd_data |
| | 21.8 MB | linear_solver_operators_m::solve_el_grp_pcg |
| | 19.0 MB | data_comm_m::copy_int_comm_to_data |
| | 18.1 MB | data_comm_m::update_int_comm |

## Peak memory

Software R&D: Observing Memory - MALT, Sébastien Valat

# Source annotations

# Tracking stack memory



Display largest stack for thread ID

Stack space used by functions on peak

Thread ID

Stack size over time

MATT WebView

Summary    Alloc sites    Time analysis    Stack    Alloc sizes    Help

Thread ID :  0

2.1 5 KB E | 54.0 KB grid_io_hdf_m::dump_grid_data_to_hdf | 29.8 KB grid_io_hdf_m::dump_grid_to_hdf | 2.1 KB | 3.6 KB | 7.9 KB solver_incompressible_

**Exascale** ∞
computing research

Example from YALES2 with gfortran issue



Many really small allocations

## Size over time

## Lifetime over size

Software R&D: Observing Memory - MALT,
Sébastien Valat

Distribution over binaries

Distribution over variables

Software R&D: Observing Memory - MALT,
Sébastien Valat

# REAL CASES

- Issue with **reallocation** on init :

- Issue only occur with **gfortran**, ifort uses stack arrays.

MATT WebView

⇅  %  |  ←  →  Allocation count ▾

Search

▮  911.9 K  data_comm_m::copy_i...

▮  896.4 K  data_comm_m::copy_...

Search intensive alloc functions

Huge number of allocation for a line programmer think it doesn't do any !

```
892    do i=1,nitem_el_grp
893        el_grp_ind = el_grp_index2int_comm_index%val(1,i)
894        int_comm_ind = el_grp_index2int_comm_index%val(2,i)
608 K  895    el_grp_r2%val(1:dim1,el_grp_ind) = int_comm_r2%val(1:dim1,int_comm_ind)
896    end do
```

**Total :**
Allocated memory : 9.5 MB
Freed memory : 9.5 MB
Max alive memory : 432
608.0 K alloc : [ 16 B , 16 B , 16 B ]
608.0 K free : [ 16 B , 16 B , 16 B ]
Lifetime : [ 24.5 K , 39.9 K , 37.8 M ] (cycles)
**Own :**
Allocated memory : 9.5 MB

And mostly really small allocations !

Software R&D: Observing Memory - MALT, Sébastien Valat

- Examples on YALES 2, small allocations :



Search for the minimal chunk size.

MATT WebView

1.0 B   /usr/lib/gcc/x86_64-p...

1.0 B   __strdup

1.0 B   data_defs_m::resize_...

Many codes produce allocations of 1B.
OK with moderation.

```
530   case (DATATYPE_REAL_NODE_VECTOR,DATATYPE_REAL_ELEM_VECTOR, &
531         DATATYPE_REAL_FACE_VECTOR,DATATYPE_REAL_PAIR_VECTOR)
532     if (associated(data_ptr%r2_ptrs)) then
533       deallocate(data_ptr%r2_ptrs)
534     end if
535     allocate(data_ptr%r2_ptrs(nel_grps))
536     do n=1,nel_grps
537       NULLIFY(data_ptr%r2_ptrs(n)%ptr)
538     end do
539
```

1 B

- Example from **Dassault mini-app** from Loïc Thébault and Eric Petit.
- **Fragmentation** can **prevent** from **returning physical pages** to OS
- **Solution** : **avoid interleaved** allocation of chunks with **different lifetime**.
- We observed with the **source annotation** that **most of them can be avoided**.

## Memory allocated over time



○ Virtual memory  ● Physical memory  ● Requested memory

3/30/2016

- On one of my own code:

```
118   //search none empty
119   bool retry;
120   do {
121       retry = false;
122       for (auto entry : channels) {
123           if (entry.acceptor()) {
124               if (entry.channel->emptyTryLock(retry) == false) {
125                   ret = entry.channel;
126                   break;
127               }
128           }
129       }
```

354.7 K (at line 122)

- It was in a **active waiting loop**….
- **Divide by 300** the number of allocation

Software R&D: Observing Memory - MALT,
Sébastien Valat

Exascale ∞
computing research



Chart — Overhead (multiplier)

| | MALT-finstr | finstr | MALT | Valgrind |
|---|---|---|---|---|
| Geant4 | 12,0 | 2,1 | 39,6 | 60,0 |
| libgeodecomp | 126,5 | 8,5 | 1,0 | 33,0 |
| libreoffice | 0,0 | 1,6 | 19,1 | 28,4 |

**Overhead (multiplier)**

Software R&D: Observing Memory - MALT,
Sébastien Valat

# Profile file size

# USAGE & CONCLUSION

- ## Backtrace mode :

  ```
  # Optionally recompile with debug flag  to get source lines :
  cc –g …
  # Run your program
  ${PREFIX}/bin/malt [--config=file.ini] YOUR_PRGM [OPTIONS]
  ```

- ## Function tracking with -finstrument-function :

  ```
  # Recompile with instrumentation flag :
  cc -finstrument-function –g …
  # Run
  ${PREFIX}/bin/malt --stack=enter-exit [--config=file.ini] YOUR_PRGM [OPTIONS]
  ```

- ## Use the web view :

  ```
  #Launch the server
  malt-webserver -i  malt-{YOUR_PRGM}-{PID}.json
  # Connect with your browser on http://localhost:8080
  ```

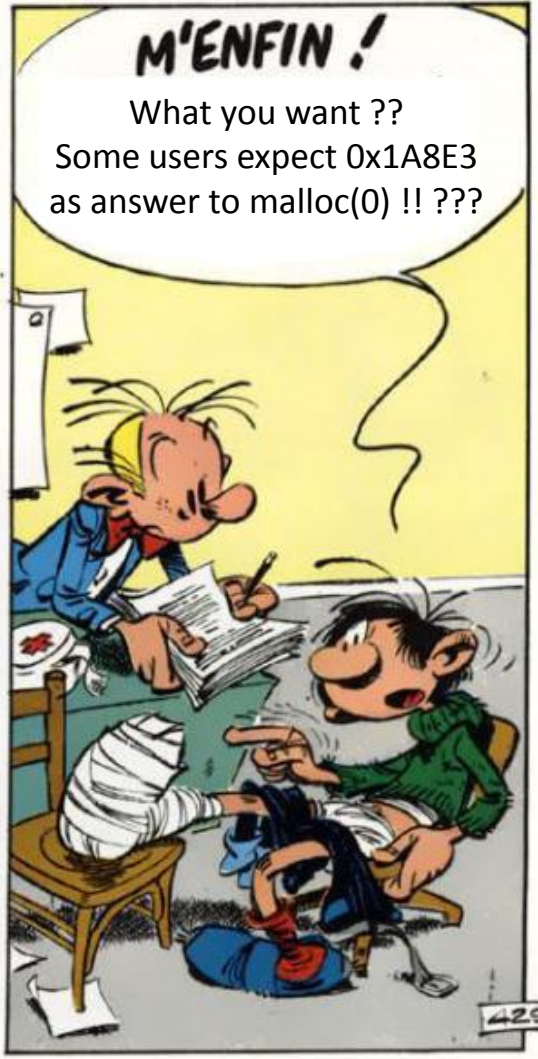Software R&D: Observing Memory - MALT,
Sébastien Valat

- Useful  tool for tracking allocation patterns

- Interesting real cases :
  - YALES2 : allocations with gfortran
  - AVBP : large allocation rate
  - PAMPA : allocation larger than expected by programmer
  - Dassault mini-app : effect and source of fragmentation

- Future work :
  - Integrate traces into the view (already get all the backend stuff)
  - Add NUMA informations (at lease statistics about mappings)
  - Hope to get Open Source release soon

Software R&D: Observing Memory - MALT,
Sébastien Valat

Thank you.

# QUESTIONS ?

# BACKUP

Software R&D: Observing Memory - MALT,
Sébastien Valat

- Add NUMA statistics

- Provide virtual/physical ratio

- Estimate page fault costs

- Exploit traces in GUI for deeper analysis
  - Alive allocations at a certain time
  - Fragmentation analysis
  - Time charts from call sites
  - Usage over threads for call sites

- The tool maintain a **call stack tree**

- Profile **stats on leafs**

- On **new** global **peak**, need **to copy** each **local current contribution**

- Need to **walk over** the wall **tree** each time ?

- Do **lazy update** :
  - Keep track of **last local peakId** on each leaf
  - On leaf update, compare the **local peakId** and the global one
  - If not same : remember the old local contribution

Software R&D: Observing Memory - MALT, Sébastien Valat